

Using High Bandwidth Memory (HBM) as an L4 Cache

- HBMs are sometimes used for very large L4 caches (128MiB to 1 GiB) instead of using HBM as a replacement for regular main memory.
- Tags and data are placed in the same HDM SDRAM row.
- First, a row access strobe (RAS) is used to get the row.
- Next, a column access strobe (CAS) is used to access the tag.
- If the tag comparison indicates a hit, then another CAS is performed to access the data, which is faster than performing another RAS.

Cache Optimization Summary

Technique	Hit time	Bandwidth	Miss penalty	Miss rate	Power consumption	Hardware cost/complexity	Comment
Small and simple caches	+			-	+	0	Trivial; widely used
Way-predicting caches	+				+	1	Used in Pentium 4
Pipelined & banked caches	-	+				1	Widely used
Nonblocking caches		+	+			3	Widely used
Critical word first and early restart			+			2	Widely used
Merging write buffer			+			1	Widely used with write through
Compiler techniques to reduce cache misses				+		0	Software is a challenge, but many compilers handle common linear algebra calculations
Hardware prefetching of instructions and data			+	+	-	2 instr., 3 data	Most provide prefetch instructions; modern high-end processors also automatically prefetch in hardware
Compiler-controlled prefetching			+	+		3	Needs nonblocking cache; possible instruction overhead; in many CPUs
HBM as additional level of cache		+/-	-	+	+	3	Depends on new packaging technology. Effects depend heavily on hit rate improvements

Where Does I/O Occur?

- between the I/O device and the cache
 - Solves the cache-coherency problem as the cache would always have the latest data.
 - Would interfere with the CPU.
- between the I/O device and main memory
 - Little interference with the CPU.
 - Can be performed in parallel with executing another process after a context switch.
 - Need to resolve the stale data problem.

Stale Data

- Copies of data can be in the data cache(s), main memory, and disk (virtual memory).
- *Stale data* occurs when these copies are not consistent.
- Two problems with stale data:
 - The CPU may not read the most recently input data since the cache may not be updated after an I/O input operation.
 - An I/O output operation may not have the correct output data because memory may not be up-to-date.

Solutions to Dealing with Stale Data

- Input
 - Not allow any I/O input buffers to be cachable.
 - Flush appropriate lines from the cache that are to be updated from the I/O input operation.
- Output
 - Use a write-through cache.
 - Write back all dirty lines in the cache before performing an I/O output operation.

Avoiding the Memory Wall

- Delays due to increasing memory latency have been largely avoided due to multiple optimizations.
 - multiple levels of cache
 - prefetching of data and instructions
 - compiler optimizations to improve locality
 - use of out-of-order pipelines with nonblocking caches supporting multiple outstanding misses
 - use of multithreading and thread-level parallelism

Fallacies and Pitfalls

- Fallacy: Predicting cache performance of one program from another.
- Pitfall: Simulating enough instructions to get accurate performance measures of the memory hierarchy.